

1 **Supplementary Information for**

2 **Unexpectedly high mutation rate of a deep-sea hyperthermophilic**
3 **anaerobic archaeon**

4
5 Jiahao Gu, Xiaojun Wang, Xiaopan Ma, Ying Sun, Xiang Xiao*, Haiwei Luo*

6
7 *Corresponding author.

8 Haiwei Luo (hluo2006@gmail.com), Xiang Xiao (zjxiao2018@sjtu.edu.cn)

9
10
11 **This PDF file includes:**

12 Supplementary Materials and Methods

13
14 References
15

16 **Materials and Methods**

17 *Sampling, cultivation, and genome sequencing of Thermococcus eurythermalis isolates*

18 Nine *Thermococcus eurythermalis* strains (Table S3) were isolated from samples of
19 Guaymas Basin hydrothermal vents in the cruise number AT 15–55, during 7-17 November
20 2009 [1]. Briefly, samples were stored in the Hungate anaerobic tubes and kept at 4°C. Then
21 the samples were enriched at 85°C or 95°C using *Thermococcales* Rich Medium (TRM). One
22 liter of TRM contains 3.3g pipes disodium salt, 23g NaCl, 5g MgCl₂·6H₂O, 0.7g KCl, 0.5g
23 (NH₄)₂SO₄, 1mL K₂HPO₄ 5%, 1mL KH₂PO₄ 5%, 1mL CaCl₂·2H₂O 2%, 0.05g NaBr, 0.01g
24 SrCl₂·6H₂O, 1mL Na₂WO₄ 10mM, 1mL FeCl₃ 25mM, 1g yeast extract, 4g tryptone and 1mg
25 resazurin [2]. The medium was adjusted to pH 7.0, autoclaved and reduced with 0.5g sodium
26 sulphide before use. Next, enrichment cultures were inoculated on the solid medium prepared
27 with hungate roll-tube technique and incubated at 85°C or 95°C under atmosphere pressure.
28 Single colonies were transferred into new TRM medium and purified using roll-tube
29 technique for 3 times and stocks were kept at -80°C. More details of sampling and isolation
30 can be found in a previous paper [1]. Among these isolates, the complete genome of the type
31 strain A501 (GCA_000769655.1) was downloaded from the NCBI GenBank database [3],
32 and the rest eight strains were sequenced in the present study. To get enrichment of these
33 eight strains, stocks kept in -80°C were inoculated into 50 mL anaerobic TRM medium in the
34 serum bottle and cultured in the incubator in 85°C. The liquid medium was supplemented
35 with sulfur and Na₂S·9H₂O. After enrichment, the cells were collected using centrifuge
36 (12,000 rpm, 10min). Genomic DNA of each isolate was extracted using the Magen Hipure

37 Soil DNA Kit and was sequenced using the Illumina Hiseq platform with 2×150 bp paired-
38 end. Raw reads were first processed by Trimmomatic 0.32 [4] to remove adaptors and trim
39 bases of low quality. The draft genome of each isolate was assembled with quality reads
40 using SPAdes v3.10.1 [5] with default parameters.

41

42 *Mutation accumulation experiment*

43 For culture propagation under high temperature, anaerobic high-temperature-tolerant
44 plates were made every day before the transfer. Plates were made using anaerobic
45 *Thermococcus* Rich Medium [2] (TRM) with gelrite (15g liter⁻¹). After sterilization, 1.5 mL
46 of a polysulfide solution [6] was added per liter of medium using syringe to make sure a
47 strictly anaerobic condition. The medium was transferred into an anaerobic chamber (COY,
48 Vinyl Anaerobic Chamber) immediately, preventing it from cooling. This is because gelrite
49 used for making plates becomes solidified soon after it become cooler. Plates were made in
50 the chamber.

51 The mutation accumulation (MA) experiment started from a single founder colony of
52 *Thermococcus eurythermalis* A501. It was transferred to new plates to form 100 independent
53 lines. Plates were put into an anaerobic jar (GeneScience), which were together moved to an
54 incubator. After incubation at 85°C under normal air pressure (optimal growth pressure from
55 0.1-30 MPa) for one day, the jar was transferred back into the anaerobic chamber. Plates were
56 then taken out. This was the initiation of the MA process. Caution was taken to ensure a
57 strictly anaerobic condition maintained throughout the experiment. A single/tiny (< 1 mm)

58 colony of each line was carefully picked and transferred onto a new plate. Then the new
59 plates were put back into the anaerobic jar for incubation. The single cell bottleneck of the
60 MA process occurred during every transfer.

61 The MA propagation was completed following 20 transfers, and four MA lines were
62 lost during the MA process. A single colony on each plate was transferred into 5 mL
63 anaerobic TRM medium in the anaerobic chamber. The liquid medium was supplemented
64 with sulfur and $\text{Na}_2\text{S}\cdot 9\text{H}_2\text{O}$. After incubation at 85°C for one day, stocks of each line were
65 kept at -80°C . Genomic DNA of each survived MA line was extracted using the Magen
66 Hipure Soil DNA Kit and sequenced using the same platform mentioned above. A
67 sequencing coverage depth of $\sim 433\times$ with an average library fragment size of ~ 470 bp was
68 obtained for each line.

69

70 *Generation time estimation with correction for cell death rate*

71 To estimate the generation time, a whole single colony was cut from 10 randomly
72 selected MA lines. The selected 10 colonies each were moved into 5 mL anaerobic TRM
73 medium supplemented with $\text{Na}_2\text{S}\cdot 9\text{H}_2\text{O}$. After dilution and re-plating, live cell density (d)
74 was measured with viable cell counts. The live and dead cell staining was done to correct the
75 total cell density for each colony. Briefly, to obtain the sufficient cell density for staining, ten
76 single colonies were cut from every MA line selected above. Live and dead bacterial staining
77 kit (Yeasen Biotech Co.) was used in this study. The kit was tested to be effective in archaea.
78 The cells were put into 350 μL anaerobic TRM medium supplemented with $\text{Na}_2\text{S}\cdot 9\text{H}_2\text{O}$.

79 After centrifuge with 10,000 g for 10 min, cells were resuspended in 50 μ L medium. Cell
80 staining was done following the protocol of the kit. Fluorescence microscope (Nikon) was
81 used to differentiate between live and dead cells. The ratio of live cells to total cells (r) was
82 0.942 (\pm 0.095) (Table S5). The number of cell divisions per transfer (D) was corrected by:

$$83 \quad D = \log_2\left(\frac{d}{r}\right)$$

84 where d is the live cell density and r is the ratio of live cells in total cells. The total number of
85 generations that each MA line went through was the multiplication of average number of cell
86 divisions per transfer and the total number of transfers. Since each MA line underwent 20
87 transfers with an average of 15.72 ± 1.76 cell divisions per transfer, there were a total of
88 314.4 ± 35.2 generations for each MA line.

89

90 *Mutation calling and mutation rate determination*

91 Raw reads were first processed by Trimmomatic 0.32 [4] to remove adaptors and trim
92 low-quality bases. Then the paired-end reads of 96 MA lines were individually mapped to the
93 *T. eurythermalis* A501 reference genome using two different mappers: BWA-mem [7] and
94 NOVOALIGN v2.08.02 (www.novocraft.com). The resulting pileup files were converted to
95 SAM format with SAMTOOLS [8].

96 The above mapping results were processed by Picard MarkDuplicates
97 (<http://broadinstitute.github.io/picard/>) to remove duplicate reads which may arise during
98 sample preparation like PCR duplication artifacts or derive from a single amplification
99 cluster. Base quality score recalibration was performed to adjust quality score affected by

100 systematic technical errors using BaseRecalibrator in GATK-4.0 [9]. Then base substitutions
101 and small indels were called using HaplotypeCaller implemented in GATK-4.0 [9]. Variants
102 were further filtered with standard parameters described by GATK Best Practices
103 recommendations, except that the Phred-scaled quality score $QUAL > 100$ and RMS
104 mapping quality $MQ > 59$ were set, which followed previous studies [9–12]. PCR primers
105 were designed with Primer Premier 5.0 [13] to confirm the presence of mutations identified
106 by the above bioinformatics method. Twenty base substitutions and nine indels were sampled
107 from 11 lines and validated. These lines were chosen because two of these lines showed the
108 highest base-substitution mutation rate and the remaining nine lines showed the longest indel
109 mutations (Table S1). The average number of analyzable sites and the average coverage per
110 site in the *T. eurythermalis* A501 MA lines were 2,123,047 (± 674) and 431 (± 57),
111 respectively.

112 The base-substitution mutation rate per nucleotide site per cell division (μ) for each
113 line was calculated according to the following equation:

$$114 \quad \mu = \frac{m}{nG}$$

115 Where m is the number of observed base substitutions, n is the number of nucleotide sites
116 analyzed, and G is the mean number of cell divisions estimated during the mutation
117 accumulation process. Following a previous study [14], the total standard error of base-
118 substitution mutation rate across all MA lines was calculated by:

$$119 \quad SE_{pooled} = \frac{s}{\sqrt{N}}$$

120 where s is the standard deviation of the mutation rate across all lines, and N is the number of
121 lines analyzed.

122

123 *The effective population size estimation for *Thermococcus eurythermalis**

124 The effective population size (N_e) of a prokaryotic species was calculated following the
125 equation $\pi_s = 2 \times N_e \times \mu$, where π_s is the nucleotide diversity at silent (synonymous) sites among
126 randomly sampled members of a species and μ is the unbiased spontaneous mutation rate.

127 Microbial species commonly harbor genetically structured populations, which has a major
128 influence on π_s and thus N_e estimation. It is therefore important to identify strains allowed for
129 free recombination when calculating N_e for a prokaryotic species [15]. The recently available
130 program PopCOGenT [16] identifies members from a prokaryotic species constituting a
131 panmictic population. The basic idea of PopCOGenT is that the recent homologous
132 recombination erased the single nucleotide polymorphisms (SNPs) and led to identical
133 regions between genomes, and therefore strains subjected with frequent recent gene transfers
134 are expected to show an enrichment of identical genomic regions compared to accumulation
135 of SNPs between genomes lacking recent transfer [16]. In practice, strains were connected
136 via recent gene flow into a network, and a putative population was identified as a cluster,
137 with within-cluster DNA transfer frequency much higher than that of between clusters. Only
138 one strain within each clonal complex was kept, which is also important for π_s estimation
139 because an overuse of strains from a clonal complex is expected to underestimate π_s . Then
140 the cluster containing the largest number of strains was chosen as the panmictic population

141 for a given species. In the case of *T. eurythermalis*, all nine strains together form a panmictic
142 population, but two strains were not used in the calculation because they were repetitive
143 members of clonal complexes.

144 Next, the single-copy orthologous genes shared by all the seven *T. eurythermalis*
145 genomes were identified by OrthoFinder 2.2.1 [17]. Amino acid sequences of each gene
146 family were aligned with MAFFT v7.464 [18] and then imposed on nucleotide sequences.
147 The number of synonymous substitution per synonymous site (d_s) for each possible gene pair
148 in each gene family was computed with the YN00 program in PAML 4.9e [19]. The π_s of
149 each gene family was obtained by averaging all pairwise d_s values, and then the median π_s
150 across all single-copy gene families together with μ were used to calculate the N_e . We used
151 the median π_s instead of the mean value, because loci showing unusually large d_s as a result
152 of allelic replacement via homologous recombination with divergent lineages are common in
153 marine prokaryotic species [20], which are expected to bias the mean value but have a limited
154 effect on the median value across gene loci. Given the small sample size of the available *T.*
155 *eurythermalis* genomes, bootstrap resampling (with replacement, 10,000 pseudoreplicates) of
156 the genomes were conducted to estimate standard deviation (Table S4) of π_s .

157

158 *Data synthesis*

159 To enable a comparative analysis of *T. eurythermalis* relative to other prokaryotic
160 species, the available μ values of other 29 prokaryotic species determined with the MA/WGS
161 technique were collected from the literature (Table S4). Among these, 20 species each had

162 multiple isolates' genomes available from the NCBI Refseq database [21], and thus were
163 used for N_e calculation. The calculation of N_e for these species followed the abovementioned
164 procedure detailed for *T. eurythermalis*, which started with the identification of members
165 constituting a panmictic population by PopCOGenT, followed by the calculation of π_S . The
166 same bootstrap resampling analysis was performed to estimate the standard deviation of π_S
167 for another four species, each of which consists of less than 10 strains (Table S4). A few
168 species have thousands of isolates' genomes available in Refseq (Table S4), which are not
169 amenable for the PopCOGenT analysis. For these species, we started from the populations
170 previously identified by ConSpeciFix [22, 23] and used these genomes as the input of
171 PopCOGenT. The ConSpeciFix delineates populations based on homoplasious SNPs, which
172 retains historical recombination signal and blurs the boundary of the ecological populations
173 enriched with recent gene transfers [16]. In the case of the species *Ruegeria pomeroyi* DSS-3,
174 a model heterotrophic marine bacterium with its mutation rate available [14], since closely
175 related isolates has not been available, we turned to its closely related species *Epibacterium*
176 *mobile* (previously known as *Ruegeria mobile*) with multiple isolates' genomes available.

177 Next, the pairwise linear relationship between μ , N_e , and genome size across the
178 prokaryotic species was initially assessed with the generalized linear model (GLM)
179 implemented in *stats* package in R v4.0.2 [24]. The Bonferonni adjusted outlier test was
180 performed with *outlierTest* function in *car* package [25]. A data point with Bonferroni p -
181 value smaller than 0.05 would be identified as the outlier. For μ versus genome size, all 30
182 species were used. In the case of N_e versus μ and N_e versus genome size, only the 21 species

183 each containing multiple strains' genomes were used. To test whether there was a
184 phylogenetic signal of these traits, the Pagel's λ [26] was estimated using the *pgls* function of
185 the *caper* package [27] which took the phylogeny of 30 species or the phylogeny of 21
186 species as an input. The species phylogeny was approximated by the 16S rRNA gene tree
187 constructed using IQ-TREE 2.0 [28] with ModelFinder [29] which assigns the best
188 substitution model and with 1,000 ultrafast bootstrap replicates. The value of λ ranges from 0
189 to 1, with 0 indicating no phylogenetic signal and 1 indicating a strong phylogenetic signal
190 due to Brownian motion. The p values for the lower and upper bounds represent whether the
191 λ is significantly different from 0 and 1, respectively. The results of this test indicate that
192 there was an intermediate phylogenetic signal for the relationship of N_e versus μ ($\lambda = 0.81$,
193 lower bound $p = 0.29$, upper bound $p = 0.06$), but not for that of N_e versus genome size and μ
194 versus genome size (in both cases, $\lambda = 0$, lower bound $p = 1$, upper bound $p < 0.001$). To
195 control for the phylogenetic effect on the correlations of the traits, the pairwise linear
196 relationship between μ , N_e , and genome size was further assessed with the phylogenetic
197 generalized least square (PGLS) regression implemented in the *caper* package [27] in R
198 v4.0.2 [24]. The PGLS and GLM regression lines were largely overlapped for N_e versus
199 genome size and μ versus genome size (Fig. 2BC). This is because no phylogenetic signal
200 was detected in these relationships. A data point was identified as an outlier in the PGLS
201 result if the associated absolute value of studentized residual is greater than three [30, 31].

202 **References**

- 203 1. Liu L, Wang F, Xu J, Xiao X. Molecular diversity of *Thermococcales* isolated from
204 Guaymas Basin hydrothermal vents. *Acta Oceanol Sin* 2013; **32**: 75–81.
- 205 2. Zeng X, Birrien J-L, Fouquet Y, Cherkashov G, Jebbar M, Querellou J, et al.
206 *Pyrococcus* CH1, an obligate piezophilic hyperthermophile: extending the upper
207 pressure-temperature limits for life. *ISME J* 2009; **3**: 873–876.
- 208 3. Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Wheeler DL. GenBank. *Nucleic*
209 *Acids Res* 2007; **35**: D21–D25.
- 210 4. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina
211 sequence data. *Bioinformatics* 2014; **30**: 2114–2120.
- 212 5. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al.
213 SPAdes: a new genome assembly algorithm and its applications to single-cell
214 sequencing. *J Comput Biol* 2012; **19**: 455–477.
- 215 6. Matsumi R, Manabe K, Fukui T, Atomi H, Imanaka T. Disruption of a sugar transporter
216 gene cluster in a hyperthermophilic archaeon using a host-marker system based on
217 antibiotic resistance. *J Bacteriol* 2007; **189**: 2683–2691.
- 218 7. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
219 transform. *Bioinformatics* 2009; **25**: 1754–1760.
- 220 8. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
221 Alignment/Map format and SAMtools. *Bioinformatics* 2009; **25**: 2078–2079.

- 222 9. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, et al. The
223 Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation
224 DNA sequencing data. *Genome Res* 2010; **20**: 1297–1303.
- 225 10. Long H, Behringer MG, Williams E, Te R, Lynch M. Similar mutation rates but highly
226 diverse mutation spectra in ascomycete and basidiomycete yeasts. *Genome Biol Evol*
227 2016; **8**: 3815–3821.
- 228 11. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A
229 framework for variation discovery and genotyping using next-generation DNA
230 sequencing data. *Nat Genet* 2011; **43**: 491–498.
- 231 12. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine
232 A, et al. From FastQ data to high confidence variant calls: the Genome Analysis Toolkit
233 best practices pipeline. *Curr Protoc Bioinformatics* 2013; **43**: 11.10.1-33.
- 234 13. Singh VK, Mangalam AK, Dwivedi S, Naik S. Primer premier: program for design of
235 degenerate primers from a protein sequence. *BioTechniques* 1998; **24**: 318–319.
- 236 14. Sun Y, Powell KE, Sung W, Lynch M, Moran MA, Luo H. Spontaneous mutations of a
237 model heterotrophic marine bacterium. *ISME J* 2017; **11**: 1713–1718.
- 238 15. Rocha EPC. Neutral theory, microbial practice: challenges in bacterial population
239 genetics. *Mol Biol Evol* 2018; **35**: 1338–1347.
- 240 16. Arevalo P, VanInsberghe D, Elsherbini J, Gore J, Polz MF. A reverse ecology approach
241 based on a biological definition of microbial populations. *Cell* 2019; **178**: 820–834.

- 242 17. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome
243 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol* 2015;
244 **16**: 157.
- 245 18. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
246 Improvements in performance and usability. *Mol Biol Evol* 2013; **30**: 772–780.
- 247 19. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 2007;
248 **24**: 1586–1591.
- 249 20. Sun Y, Luo H. Homologous recombination in core genomes facilitates marine bacterial
250 adaptation. *Appl Environ Microbiol* 2018; **84**: e02545-17.
- 251 21. O’Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, et al. Reference
252 sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and
253 functional annotation. *Nucleic Acids Res* 2016; **44**: D733-745.
- 254 22. Bobay L-M, Ellis BS-H, Ochman H. ConSpeciFix: classifying prokaryotic species
255 based on gene flow. *Bioinformatics* 2018; **34**: 3738–3740.
- 256 23. Bobay L-M, Ochman H. Biological species are universal across life’s domains. *Genome*
257 *Biol Evol* 2017; **9**: 491–501.
- 258 24. R Development Core Team. R: A language and environment for statistical computing. R
259 Foundation for Statistical Computing, Vienna, Vienna, Austria.
- 260 25. Fox J, Weisberg S. An R Companion to Applied Regression, Third Edition. 2019. Sage,
261 Thousand Oaks CA.

- 262 26. Pagel M. Inferring the historical patterns of biological evolution. *Nature* 1999; **401**:
263 877–884.
- 264 27. Orme D. The caper package: comparative analysis of phylogenetics and evolution in R.
265 2013; 36.
- 266 28. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective
267 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol*
268 2015; **32**: 268–274.
- 269 29. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermin LS. ModelFinder:
270 fast model selection for accurate phylogenetic estimates. *Nat Methods* 2017; **14**: 587–
271 589.
- 272 30. Powell LE, Isler K, Barton RA. Re-evaluating the link between brain size and
273 behavioural ecology in primates. *Philos Trans R Soc Lond, B, Biol Sci* 2017; **284**:
274 20171765.
- 275 31. Harrison TL, Simonsen AK, Stinchcombe JR, Frederickson ME. More partners, more
276 ranges: generalist legumes spread more easily around the globe. *Biol Lett* 2018; **14**:
277 20180616.
278